

# Satisfiability of a Spatial Logic with Tree Variables

Emmanuel Filiot<sup>1</sup>   Jean-Marc Talbot<sup>2</sup>   Sophie Tison<sup>1</sup>

<sup>1</sup> INRIA Futurs, Mostrare Project, University of Lille 1 (LIFL, UMR 8022 of CNRS)

<sup>2</sup> University of Provence (LIF, UMR 6166 of CNRS), Marseille

**Abstract.** We investigate in this paper the spatial logic TQL for querying semi-structured data, represented as unranked ordered trees over an infinite alphabet. This logic consists of usual Boolean connectives, spatial connectives (derived from the constructors of a tree algebra), tree variables and a fixpoint operator for recursion. Motivated by XML-oriented tasks, we investigate the guarded TQL fragment. We prove that for closed formulas this fragment is MSO-complete. In presence of tree variables, this fragment is strictly more expressive than MSO as it allows for tree (dis)equality tests, i.e. testing whether two subtrees are isomorphic or not. We devise a new class of tree automata, called TAGED, which extends tree automata with global equality and disequality constraints. We show that the satisfiability problem for guarded TQL formulas reduces to emptiness of TAGED. Then, we focus on bounded TQL formulas: intuitively, a formula is bounded if for any tree, the number of its positions where a subtree is captured by a variable is bounded. We prove this fragment to correspond with a subclass of TAGED, called bounded TAGED, for which we prove emptiness to be decidable. This implies the decidability of the bounded guarded TQL fragment. Finally, we compare bounded TAGED to a fragment of MSO extended with subtree isomorphism tests.

## 1 Introduction

In this paper, we consider the spatial logic TQL [7]. Spatial logics have been used to express properties of structures such as trees [7], graphs [6, 12] and heaps [19]. The main ingredients of spatial logics are spatial connectives: roughly speaking, these connectives are derived from operators that can be used to generate the domain of interpretation.

The logic we consider here is interpreted over hedges (*i.e.* unranked ordered trees) labelled over an infinite alphabet. The logic integrates Boolean connectives, spatial connectives (derived from the constructors of an unranked ordered tree algebra), tree variables and a fixpoint operator for recursion.

We focus on the satisfiability problem of TQL. It is quite simple to prove that the full TQL logic over unranked ordered trees even without tree variables is undecidable. We then focus on the guarded fragment which ensures that recursive variables have to be guarded by tree extension. We show that guarded TQL logic without tree variables is equivalent to the monadic second order logic (MSO).

However, when tree variables are considered, things are getting more complicated. Indeed, we can express that two non-empty paths starting from a node of a tree lead to two isomorphic subtrees, which goes beyond regularity over unranked trees.

Still about expressiveness of this logic, an infinite alphabet and the ability to test for tree equality allow us to consider some data values. We can write formulas whose models are hedges which violate some key constraints or some functional dependencies.

We focus on bounded TQL formulas: intuitively, a formula is bounded if for any tree, the number of equalities and disequalities that have to be tested – to check non-deterministically that the tree is a model of the formula – is bounded.

We introduce a new class of tree automata, called *tree automata with global equalities and disequalities* (TAGED for short), which extends unranked tree automata  $A$  with an equality relation  $=_A$  and a disequality relation  $\neq_A$  on states. Subtrees of some tree  $t$  which evaluates by  $A$  to states which are in relation by  $=_A$  (resp. by  $\neq_A$ ) have to be isomorphic (resp. non isomorphic). Naturally,  $=_A$  induces an equivalence relation on a subset of nodes of  $t$ , but the number of classes of this relation is bounded. E.g., TAGED can express that all subtrees of height  $n$ , for some fixed natural  $n$ , are equal, but not that for each node of the tree, all the subtrees rooted at its sons are equal. Although it is a natural extension of tree automata, this extension has never been considered.

We show that satisfiability of guarded TQL formulas reduces to emptiness of TAGED. We define a subclass of TAGED, called *bounded TAGED*, for which we can decide emptiness. Intuitively, boundedness ensures that the cardinality of every equivalence class is bounded, which may not be the case for full TAGED. We show emptiness decidability of bounded TAGED.

We complete our proof by constructing a TAGED from a guarded and bounded TQL formula. This construction extends the one from [4] with tree variables. This extension is non-trivial as the automata we have to consider are non-deterministic.

Finally, we define an extension of MSO with a binary relation  $\sim$  between nodes; two nodes are in relation if they are roots of two isomorphic subtrees. We consider MSO formulas extended with the predicate  $\sim$ . It is easy to see that this extension renders MSO undecidable. However, we prove that if the relation  $\sim$  concerns only variables belonging to a prefix of existentially quantified first-order variables, then this extension is decidable. The proof works by reduction to emptiness of bounded TAGED.

Automata dealing with data values have been studied in [18, 3]. However, our motivations are different and we obtain the capability to manage data values as a side-effect. In [3], the authors study two-variables FO logic extended with an equality relation on data values. The expressiveness of this formalism and the one presented here are not comparable: we can test tree isomorphisms while they can test data value equality only, but restricting our logic to data-value equality is strictly less expressive, as we do not have quantifiers.

The paper is organised as follows: in Section 2 we recall definitions for hedges, hedge automata and monadic second order logic. Section 3 describes the TQL logic and results we obtain concerning its satisfiability. Section 4 is dedicated to the tool we use to solve the satisfiability problem, namely tree automata with global equalities and disequalities (TAGED). In Section 5 we relate satisfiability of TQL formulas and emptiness of TAGED. Finally, in Section 6 we propose an extension of MSO with isomorphism tests whose satisfiability problem is decidable.

## 2 Preliminaries

We consider an infinite set of labels  $A$ .

*Hedges - Trees* Let  $\Sigma$  be the signature  $\{0, |\} \cup \{a \mid a \in \Lambda\}$ , where 0 is a constant, | a binary symbol and  $a$ s unary symbols. We call *hedge* an element of the  $\Sigma$ -algebra *Hedge* obtained by quotienting the free  $\Sigma$ -algebra by the following three axioms:

$$0|h = h \quad h|0 = h \quad (h|h'')|h''' = h|(h''|h''')$$

0 will be the *empty hedge*. We call respectively *trees* and *leaves* hedges of the form  $a(h)$  and  $a(0)$ . We may omit | and write  $a(h)b(h')c(h'')$  instead of  $a(h)|b(h')|c(h'')$ . We define  $\text{roots}(h)$  as the word from  $\Lambda^*$  defined recursively as : (i)  $\text{roots}(0) = \epsilon$ , (ii)  $\text{roots}(a(h')) = a$  and, (iii)  $\text{roots}(h_1|h_2) = \text{roots}(h_1)\text{roots}(h_2)$ .

We will also adopt the graph point of view and consider hedges as a set of vertices  $V$ , two disjoint sets of directed edges  $E_c, E_s$  and a mapping  $\lambda$  from  $V$  to  $\Lambda$ . In a hedge  $h$ , one associates a vertex with each occurrence of elements of  $\Lambda$ . There is an edge from  $E_c$  (resp. from  $E_s$ ) from an occurrence  $a_1$  to an occurrence  $a_2$  if the hedge contains a pattern of the form  $a_1(h_1|a_2(h)|h_2)$  for some hedges  $h_1, h_2$  (resp.  $a_1(h_1)|a_2(h_2)$  for some hedges  $h_1, h_2$ ).

For every hedge  $h = (V, E_c, E_s, \lambda)$ , we denote by  $\text{nodes}(h)$  the set  $V$  and by  $\text{lab}_h(u)$  the label  $\lambda(u)$ , for  $u \in V$ . We denote  $h|_u$  the subtree of  $h$  rooted at  $u$ , and by  $\leq$  the reflexive-transitive closure of  $E_c$ . E.g., the root is minimal for  $\leq$  in a tree.

For a set of labels  $L$ , we denote  $\mathbb{H}_L$  (resp.  $\mathbb{T}_L$ ) the set of hedges (resp. trees) with nodes labelled by elements in  $L$ .

*Hedge automata [17]* A *hedge automaton*  $A$  is a 4-tuple  $(\Lambda, Q, F, \Delta)$  where  $\Delta$  is a finite set of rules  $\alpha(L) \rightarrow q$  where  $\alpha$  is a finite or cofinite set of labels,  $L \subseteq Q^*$  is a regular language over states from  $Q$ , and  $F \subseteq Q^*$  is an accepting regular language.

**Definition 1 (runs).** Let  $h$  be a hedge and  $A$  be a hedge automaton. The set of runs  $R_A(h) \subseteq \mathbb{H}_Q$  of  $A$  on  $h$  is the set of hedges over  $Q$  inductively defined by:

$$\begin{aligned} R_A(h_1|h_2) &= \{r_1|r_2 \mid r_1 \in R_A(h_1), r_2 \in R_A(h_2)\} & R_A(0) &= \{0\} \\ R_A(a(h)) &= \{q(r) \mid \exists \alpha(L) \rightarrow q \in \Delta, a \in \alpha, r \in R_A(h), \text{roots}(r) \in L\} \end{aligned}$$

Let  $\bar{q}$  be a word of states, we denote by  $R_{A, \bar{q}}(h) \subseteq R_A(h)$  the set of runs  $r$  such that  $\text{roots}(r) = \bar{q}$ , and often say that  $h$  evaluates to  $\bar{q}$  by  $A$ . The set of accepting runs of  $A$  on  $h$ , denoted by  $R_A^{\text{acc}}(h)$ , is defined by  $\{r \mid \exists \bar{q} \in F, r \in R_{A, \bar{q}}(h)\}$ .

The language accepted by  $A$ , denoted  $L(A)$ , is defined by  $\{h \mid R_A^{\text{acc}}(h) \neq \emptyset\}$ .

Testing emptiness of the language accepted by a hedge automaton is decidable [5].

*MSO* The logic MSO (Monadic Second Order logic) is the extension of the first-order logic FO with the possibility to quantify over unary relations, i.e. over sets.

Let  $\sigma$  be the signature  $\{\text{lab}_a \mid a \in \Lambda\}$  where  $\text{lab}_a$ s are unary predicates. We associate with a hedge  $h = (V, E_c, E_s, \lambda)$  a finite  $\sigma$ -structure  $\mathcal{S}^h = \langle V, \{\text{ch}, \text{ns}\} \cup \{\text{lab}_a^h \mid a \in \Lambda\} \rangle$ , such that  $\text{lab}_a^h(v)$  (resp.  $\text{ch}(v, v')$ ,  $\text{ns}(v, v')$ ) holds in  $\mathcal{S}^h$  if  $\lambda(v) = a$  (resp.  $(v, v') \in E_c$ ,  $(v, v') \in E_s$ ).

We assume a countable set of first-order variables ranging over by  $x, y, z, \dots$  and a countable set of second-order variables ranging over by  $X, Y, Z, \dots$

MSO formulas are given by the following syntax:

$$\psi ::= \text{lab}_a(x) \mid \text{ch}(x, y) \mid \text{ns}(x, y) \mid x \in X \mid \psi \vee \psi \mid \neg \psi \mid \exists x. \psi \mid \exists X. \psi$$

$\phi ::= 0$	empty hedge	$\llbracket 0 \rrbracket_{\rho, \delta}$	$= \{0\}$
$\top$	truth	$\llbracket \top \rrbracket_{\rho, \delta}$	$= \mathbb{H}_A$
$\alpha[\phi]$	extension	$\llbracket \alpha[\phi] \rrbracket_{\rho, \delta}$	$= \{a(h) \mid h \in \llbracket \phi \rrbracket_{\rho, \delta}, a \in \alpha\}$
$\phi \phi$	composition	$\llbracket \phi \phi' \rrbracket_{\rho, \delta}$	$= \{h h' \mid h \in \llbracket \phi \rrbracket_{\rho, \delta}, h' \in \llbracket \phi' \rrbracket_{\rho, \delta}\}$
$\neg\phi$	negation	$\llbracket \neg\phi \rrbracket_{\rho, \delta}$	$= \mathbb{H}_A \setminus \llbracket \phi \rrbracket_{\rho, \delta}$
$\phi \vee \phi$	disjunction	$\llbracket \phi \vee \phi' \rrbracket_{\rho, \delta}$	$= \llbracket \phi \rrbracket_{\rho, \delta} \cup \llbracket \phi' \rrbracket_{\rho, \delta}$
$X$	tree variable	$\llbracket X \rrbracket_{\rho, \delta}$	$= \{\rho(X)\}$
$\xi$	recursion variables	$\llbracket \xi \rrbracket_{\rho, \delta}$	$= \delta(\xi)$
$\mu\xi.\phi$	least fixpoint	$\llbracket \mu\xi.\phi \rrbracket_{\rho, \delta}$	$= \bigcap \{S \subseteq \mathbb{H}_A \mid \llbracket \phi \rrbracket_{\rho, \delta[\xi \mapsto S]} \subseteq S\}$
$\phi^*$	iteration	$\llbracket \phi^* \rrbracket_{\rho, \delta}$	$= 0 \cup \bigcup_{i>0} \underbrace{\llbracket \phi \rrbracket_{\rho, \delta} \mid \dots \mid \llbracket \phi \rrbracket_{\rho, \delta}}_{i \text{ times}}$

(a) Syntax
(b) Semantics

**Fig. 1.** TQL logic

Let  $S$  be a  $\sigma$ -structure with domain  $V$ . Let  $\rho$  be a valuation mapping first-order variables to elements from  $V$  and second-order variables to subsets of  $V$ . We write  $S \models_\rho \psi$  when the structure  $S$  is a model of the formula  $\psi$  under the valuation  $\rho$ ; this is defined in the usual Tarskian manner and we have in particular, (i)  $\psi \models_\rho \text{lab}_a(x)$  if  $\text{lab}_a(\rho(x))$  holds in  $S$ , (ii)  $\psi \models_\rho \text{ch}(x, y)$  if  $\text{ch}(\rho(x), \rho(y))$  holds in  $S$ , (iii)  $\psi \models_\rho \text{ns}(x, y)$  if  $\text{ns}(\rho(x), \rho(y))$  holds in  $S$ .

A set of hedges  $S$  is *MSO-definable* if there is an MSO sentence  $\psi$  such that  $S = \{h \mid h \models \psi\}$ . It is well-known that a language is accepted by some hedge automata iff it is MSO-definable.

### 3 The Tree Query Logic

We consider here a fragment of the TQL logic defined in [7] and adapt it to unranked ordered trees.

*Syntax* We assume a countable set  $\mathcal{T}$  of tree variables ranging over by  $X, Y$ , and a countable set  $\mathcal{R}$  of recursion variables ranging over by  $\xi$ . Let  $\alpha$  be a finite or co-finite set of labels from  $A$ . Formulas  $\phi$  from TQL are given by the syntax on Figure 1(a). We allow cofinite sets in extensions, otherwise we could not express formula  $A[0]$ .

We assume that  $\mu$  is the binder for recursion variables and the notions of bound and free variables are defined as usual. To ensure the existence of fixpoint, we will assume that in formulas  $\mu\xi.\phi$ , the recursion variable  $\xi$  occurs under an even number of negations. A formula is said to be *recursion-closed* if all the occurrences of its recursion variables are bound. A TQL *sentence* is a recursion-closed formula that does not contain tree variables. A TQL formula  $\phi$  is *guarded* if for any of its subformula  $\mu\xi.\phi'$ , the variable  $\xi$  occurs in the scope of some extension operator  $\alpha[\ ]$  in  $\phi'$ .

We assume from now on that recursion variables are bound only once in formulas and denote  $\text{recvar}(\phi)$  (resp.  $\text{var}(\phi)$ ) the set of recursion variables (resp. tree variables) occurring in  $\phi$ . We may write  $a[\phi]$  instead of  $\{a\}[\phi]$ .

*Semantics* Interpretation maps a TQL formula to a set of hedges. Let  $\rho$  be an assignment of tree variables into  $\mathbb{T}_A$  and  $\delta$  be an assignment of the recursion variables into sets of hedges. The interpretation of the formula  $\phi$  under  $\rho$  and  $\delta$ , denoted by  $\llbracket \phi \rrbracket_{\rho, \delta}$  is inductively defined and given on Figure 1(b).

*Examples* Let us consider the following formulas:

$$\phi = a[\top]|\top \quad (1)$$

$$\phi_s = \mu_{\xi}.(a[\top]|\xi \vee 0) \quad (2)$$

$$\phi_{dtd} = (employee[name[\Lambda[0]] \mid dpt[\Lambda[0]] \mid manager[\Lambda[0]])^* \quad (3)$$

$$\phi_{dd} = \phi_{dtd} \wedge \top | employee[X] | \top | employee[X] | \top \quad (4)$$

$$\phi_{path(a),0} = \mu\xi.((\top|a[\xi]|\top) \vee 0) \quad (5)$$

The above formula  $\phi$  is interpreted as the set of hedges of length at least 1, such that the root of the first tree is labelled  $a$ . The formula  $\phi_s$  is interpreted as the set of hedges  $\{a[h_1] \dots [a[h_n]] \mid h_i \in \mathbb{H}_A, n \geq 0\}$ . The formula  $\phi_{dtd}$  is interpreted as the set of hedges defining employees by their name, the department they work in, and their manager whereas the formula  $\phi_{dd}$  expresses that an employee occurs twice in the database. Finally, the models of the formula  $\phi_{path(a),0}$  are hedges with a path labelled by  $as$  from one of the roots to some leaf (*i.e.* the empty hedge 0).

The formula  $\phi_{odd}$  is interpreted as the set of hedges having an odd number of nodes:

$$\phi_{odd} = \mu\xi_o.(\Lambda[0] \vee \Lambda[\phi_{even}(\xi_o)]\phi_{even}(\xi_o) \vee \Lambda[\xi_o]|\xi_o$$

$$\text{where } \phi_{even}(\xi_o) = \mu\xi_e.(0 \vee (\Lambda[\xi_o]|\xi_e \vee \Lambda[\xi_e]|\xi_o))$$

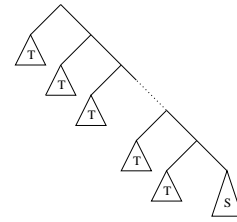
Let us denote  $\phi_{path(L),\psi} = \mu\xi.((\top | L[\xi] | \top) \vee \psi)$  the formula whose models are hedges containing a path labelled by elements from  $L$  from one of the roots to a hedge satisfying  $\psi$ . The models of the following formula are the hedges having two non-empty paths labelled respectively by  $as$  and by  $bs$  from two of the top-level nodes, those two paths leading to some identical subtrees

$$\top \mid (a[\phi_{path(a),X}] \mid \top \mid b[\phi_{path(b),X}] \vee b[\phi_{path(b),X}] \mid \top \mid a[\phi_{path(a),X}] \mid \top$$

The formula  $\phi_{id\_not\_key}$  is interpreted as the set of hedges for which two nodes labelled  $id$  have identical subtrees (roughly speaking “the (data) value of the element  $id$  can not be used as a key in this XML document”)

$$\phi_{id\_not\_key} = \top | \Lambda[\phi_{path(\Lambda), id[X]}] | \top | \Lambda[\phi_{path(\Lambda), id[X]}] | \top$$

The following formula states that two trees *employee* have identical subtrees rooted by *dpt* but different subtrees rooted by *manager*



**Fig.2.** A tree with  $T \neq S$

$$\phi_{dtd} \wedge \top \mid \text{employee}[\top \mid \text{dpt}[X] \mid \text{manager}[Y]] \mid \top$$

$$\mid \text{employee}[\top \mid \text{dpt}[X] \mid \text{manager}[\neg Y]] \mid \top$$

A hedge satisfying this formula may be considered as ill-formed assuming the existence of some functional dependency stating that *department* has only one *manager*.

Models of the formula  $\phi_{branch}$  are the trees whose shapes are described on Figure 2.

$$\phi_{branch} = a[X|\mu\xi.(a[X|\xi] \vee \neg X \wedge \Lambda[\top])]$$

$\phi_s$  and  $\phi_{odd}$  are the only two formulas that are not guarded; however,  $\phi_s$  is equivalent to  $a[\top]^*$ , which is guarded and  $\phi_{odd}$  to the following guarded formula

$$\begin{aligned}\phi_{odd} &= \mu\xi_o. \Lambda[\phi_{even}(\xi_o) \vee (\xi_o|\phi_{even}(\xi_o)^*|\xi_o)^*] \\ \phi_{even}(\xi_o) &= \mu\xi_e. \Lambda[\xi_e^*\xi_o|(\phi_{even}(\xi_o) \vee (\xi_o|\phi_{even}(\xi_o)^*|\xi_o)^*) \vee 0\end{aligned}$$

But the formula  $\mu\xi.(a[0]|\xi|b[0] \vee 0)$  is neither guarded nor equivalent to any guarded formula.

**Definition 2 (satisfiability).** A recursion-closed TQL formula  $\phi$  is satisfiable if there exists a hedge  $h$  and an assignment  $\rho$  such that  $h \in \llbracket \phi \rrbracket_\rho$ .

*TQL sentences* It is easy to prove that TQL sentences can describe sets of hedges that are not MSO-definable; for instance, the TQL sentence  $\mu\xi.(a[0]|\xi|b[0] \vee 0)$  describes a “flat” hedge of the form  $(a[0]^n b[0]^n)$  for  $n \in \mathbb{N}$ .

**Theorem 1.** For any set of hedges  $S$ , there exists a guarded TQL sentence  $\phi$  such that  $\llbracket \phi \rrbracket = S$  iff  $S$  is MSO-definable.

As a consequence of Theorem 4, satisfiability is decidable for guarded TQL sentences. This restriction is crucial, since, by reduction of emptiness problem for the intersection of two context-free grammars, one can prove that:

**Theorem 2.** Satisfiability for TQL sentences is undecidable.

*Adding quantification* As in [7], one could also consider quantification over tree variables  $\exists X$  and  $\forall X$  with the following semantics:

$$\llbracket \exists X.\phi \rrbracket_{\rho,\delta} = \bigcup_{t \in \mathbb{T}} \llbracket \phi \rrbracket_{\rho[X \rightarrow t],\delta} \quad \llbracket \forall X.\phi \rrbracket_{\rho,\delta} = \bigcap_{t \in \mathbb{T}} \llbracket \phi \rrbracket_{\rho[X \rightarrow t],\delta}$$

where  $\rho[X \rightarrow t]$  is the assignment identical to  $\rho$  except that it associates  $t$  with  $X$ .

Hence, the satisfiability problem from Definition 2 is equivalent to the one of closed formulas of the form  $\exists X_1 \dots \exists X_n \phi$  where  $\phi$  is a recursion-closed TQL formula, i.e. the existence of a tree satisfying this formula.

For more complicated alternation of quantifiers, one can easily adapt the proof from [8] about the undecidability of the fragment of TQL without iteration, recursion and tree variable but with quantification over labels to prove that

**Theorem 3.** Satisfiability for recursion-closed TQL formulas with quantification is undecidable (this holds even for recursion-free formulas).

*Bounded TQL formulas* In bounded formulas, variables can occur in recursion only in a restricted manner: intuitively a formula is bounded if there exists a bound on the number of equality test performed to (non-deterministically) verify that a hedge is a model of the formula. Boundedness appears naturally in unification problems and in pattern languages (where variables appears a bounded number of times in terms or patterns). But defining boundedness in the presence of recursion is a bit more technical.

In the examples we have given so far, the only formula that is not bounded is  $\phi_{branch}$ . The following formula is also not bounded as it asks each subtree of the hedge to be different from  $X$ :  $\neg(\mu\xi.\top|(X \vee \Lambda[\xi])|\top)$

We let  $\beta$  be a recursion-closed formula s.t. no recursion variable is bound twice, and denote by  $\phi_\xi$  the formula s.t.  $\mu\xi.\phi_\xi$  is a subformula of  $\beta$ . To define *boundedness* formally, we introduce, for every subformula  $\phi$  of  $\beta$ , the *generalized free variables* of  $\phi$ , denoted  $\text{var}_\beta^*(\phi)$ , as the least solution of the following (recursive) equations:

$$\begin{aligned} \text{var}_\beta^*(0) &= \text{var}_\beta^*(\top) = \emptyset & \text{var}_\beta^*(a[\phi]) &= \text{var}_\beta^*(\neg\phi) = \text{var}_\beta^*(\mu\xi.\phi) = \text{var}_\beta^*(\phi^*) = \text{var}_\beta^*(\phi) \\ \text{var}_\beta^*(\xi) &= \text{var}_\beta^*(\phi_\xi) & \text{var}_\beta^*(X) &= \{X\} & \text{var}_\beta^*(\phi \vee \phi') &= \text{var}_\beta^*(\phi|\phi') = \text{var}_\beta^*(\phi) \cup \text{var}_\beta^*(\phi') \end{aligned}$$

Note that the least solution of these equations is computable. An operator occurs positively (resp. negatively) in a formula if it occurs under an even (resp. odd) number of negations.

A formula  $\beta$  is *bounded* if it satisfies the following properties:

1. for any subformula  $\phi^*$ ,  $\text{var}_\beta^*(\phi) = \emptyset$ ;
2. for any subformula  $\phi|\phi'$  where  $|$  occurs negatively,  $\text{var}_\beta^*(\phi) = \text{var}_\beta^*(\phi') = \emptyset$ .
3. each formula  $\phi|\phi'$  where  $|$  occurs positively and each formula  $\phi \vee \phi'$  where  $\vee$  occurs negatively satisfy  $\forall \xi \in \text{recvar}(\phi), \text{var}_\beta^*(\xi) \cap \text{var}_\beta^*(\phi') = \emptyset$ , and  $\forall \xi' \in \text{recvar}(\phi'), \text{var}_\beta^*(\xi') \cap \text{var}_\beta^*(\phi) = \emptyset$ .

As a consequence of Theorem 1, this fragment is strictly more expressive than guarded TQL sentences, as it allows for tree isomorphism tests. Combining Theorems 6 and 5 of the next two sections proves our main result:

**Theorem 4.** *Satisfiability is decidable for bounded guarded TQL formulas.*

*Remarks* Our logic can be seen as an extension of the (recursive) pattern-language of XDuce [13]. The main difference here is that we allow Boolean operators and drop the *linear* condition for variables of XDuce. The pattern-matching mechanism of CDuce [1] extends the one from XDuce with Boolean operations and weaker conditions on variables. However, no equality tests between terms can be performed making our logic more powerful. Since we consider an infinite alphabet and we allow equality tests between trees, we can, as a side effect, simulate data values as done in some of the examples we gave.

## 4 Tree Automata with Global Equalities and Disequalities

In this section, we present a new extension of hedge automata (called TAGED) with the ability to test tree equalities or disequalities globally on the run. We prove decidability of emptiness for a particular class of TAGED, called bounded TAGED, which we use to decide satisfiability of bounded TQL formulas.

### 4.1 Definitions

**Definition 3 (TAGED).** *A tree automaton with global equalities and disequalities (TAGED) is a 6-tuple  $A = (\Lambda, Q, F, \Delta, =_A, \neq_A)$  such that:*

- $(\Lambda, Q, F, \Delta)$  is a hedge automaton;
- $=_A$  is an equivalence relation on a subset of  $Q$ ;
- $\neq_A$  is a non-reflexive symmetric binary relation on  $Q$ ;

A TAGED is *positive* if  $\neq_A = \emptyset$ , and is simply denoted by  $A = (\Lambda, Q, F, \Delta, =_A)$ . The set  $\{q \mid \exists q' \in Q, q =_A q'\}$  is denoted by  $E$ , and for all states  $q \in E$ , we denote by  $[q]$  its equivalence class. The set  $\{q \mid \exists q' \in Q, q \neq_A q'\}$  is denoted by  $D$ . The notion of run differs from hedge automata as we add equality and disequality constraints.

**Definition 4 (runs).** Let  $A = (\Lambda, Q, F, \Delta, =_A, \neq_A)$  be a TAGED. A run  $r$  of the hedge automaton  $(\Lambda, Q, F, \Delta)$  on a hedge  $h$  satisfies the equality constraints if the following holds:  $\forall i \in \{1, \dots, n\}, \forall u, v \in \text{nodes}(h), \text{lab}_r(u) =_A \text{lab}_r(v) \implies h|_u = h|_v$ .

Similarly, the run  $r$  on  $h$  satisfies the disequality constraints if the following holds:  $\forall i \in \{1, \dots, n\}, \forall u, v \in \text{nodes}(h), \text{lab}_r(u) \neq_A \text{lab}_r(v) \implies h|_u \neq h|_v$ .

The set of accepting runs of  $A$  on  $h$ , denoted  $R_A^{\text{acc}}(h)$ , is the set of accepting runs  $r$  of  $(\Lambda, Q, F, \Delta)$  which satisfy the equality and disequality constraints. The language accepted by  $A$ , denoted  $L(A)$ , is the set of hedges  $h$  such that  $R_A^{\text{acc}}(h) \neq \emptyset$ .

Remark that  $L(A)$  is not necessarily regular, as illustrated by Example 1.

*Example 1.* Let  $\Lambda$  be an infinite alphabet. Let  $Q = \{q, q_X, q_f\}$ ,  $F = \{q_f\}$ , and let  $\Delta$  be defined as the set of following rules:  $\Lambda(q^*) \rightarrow q \quad \Lambda(q^*) \rightarrow q_X \quad a(q_X q_X) \rightarrow q_f$ . Let  $A_1$  be the positive TAGED  $(\Lambda, Q, F, \Delta, \{q_X =_{A_1} q_X\})$ . Then  $L(A_1)$  is the set  $\{a(t|t) \mid a \in \Lambda, t \in \mathbb{T}_\Lambda\}$ , which is known to be non regular [10].

*Example 2.* Let  $Q = \{q, q_X, q_{\overline{X}}, q_f\}$ ,  $F = \{q_f\}$ , and let  $\Delta$  defined as the set of following rules:

$$\Lambda(q^*) \rightarrow q_{\overline{X}} \quad \Lambda(q^*) \rightarrow q \quad \Lambda(q^*) \rightarrow q_X \quad a(q_X(q_{\overline{X}} + q_f)) \rightarrow q_f$$

Let  $A_2$  be the TAGED  $(\Lambda, Q, F, \Delta, \{q_X =_{A_2} q_X\}, \{q_X \neq_{A_2} q_{\overline{X}}\})$ . Then  $L(A_2)$  is the set of hedges whose shapes are described on Figure 2.

*Remarks* Extensions of tree automata which allow for syntactic equality and disequality tests between subterms have already been defined by adding constraints to automaton rules. E.g., adding the constraint  $1.2 = 2$  to a rule means that one can apply the rule at position  $\pi$  only if the subterm at position  $\pi.1.2$  is equal to the subterm at position  $\pi.2$ . Testing emptiness of the recognized language is undecidable in general [16] but two classes with a decidable emptiness problem have been emphasised. In the first class, automata are deterministic and the number of equality tests along a path is bounded [11] whereas the second restricts tests to sibling subterms [2]. This latter class has recently been extended to unranked trees [15], the former one has been extended to equality modulo equational theories [14]. But, contrarily to TAGED, in all these classes, tests are performed locally, typically between sibling or cousin subterms. Finally, automata with local and global equality tests, using one memory, have been considered in [9]. Their emptiness problem is decidable, and they can simulate positive TAGEDs which use at most one state per runs to test equalities, but not general positive TAGEDs.

**Definition 5 (bounded TAGED).** A bounded TAGED is a 7-tuple  $A = (\Lambda, Q, F, \Delta, =_A, \neq_A, k)$  where  $A' = (\Lambda, Q, F, \Delta, =_A, \neq_A)$  is a TAGED and  $k \in \mathbb{N}$  is a natural. An accepting run  $r$  of  $A$  on a tree  $t$  is an accepting run of  $A'$  on  $t$  such that the following is true:  $|\{u \mid \text{lab}_r(u) \in E \cup D\}| \leq k$ .

We say that  $A$  and its accepting runs are  $k$ -bounded and may write  $(A', k)$  instead of  $A$ . We say that a TAGED  $B$  is equivalent to a bounded TAGED  $A$  if  $L(A) = L(B)$ .



The TAGED of Example 1 is equivalent to the 2-bounded TAGED  $(A_1, 2)$ , whereas one can prove that the one of Example 2 is not equivalent to any bounded TAGED.

**Theorem 5 (emptiness of bounded TAGED).** *Let  $A$  be a bounded TAGED. It is decidable to know whether  $L(A) = \emptyset$  holds or not.*

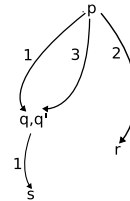
The rest of this subsection is dedicated to the proof of this theorem, first for positive bounded TAGED, then for full bounded TAGED.

## 4.2 Configurations

We define a tool called *configurations* used to decide emptiness of positive bounded TAGED. In this subsection, the 6-tuple  $A = (\Lambda, Q, F, \Delta, =_A, k)$  always denotes a positive bounded TAGED. We suppose that  $A$  accepts trees only, i.e.  $F \subseteq Q$ . It is not difficult to adapt the decidability result to hedge acceptors. Moreover, we suppose that for any run  $r$  -even non accepting- on a tree, the cardinal of the set of nodes labelled by states of  $E$  is at most  $k$ . Indeed, it is easy to transform  $A$  to ensure this property by enriching states with a counter up to  $k$ . We show how to decompose a positive TAGED into an equivalent and computable finite set of configurations. Since testing emptiness of configurations is easily decidable, we get the decidability result for positive TAGED. Informally, configurations are (non-regular) tree acceptors which make explicit parent or ancestor relations between nodes for which equality tests are performed by  $A$ . These are DAG-like structures labelled by sets of states of  $A$ . A tree  $t$  is accepted by some configuration  $c$  if the unfolding of  $c$  can be embedded into a run  $r$  of  $A$  on  $t$ , such that labels of  $r$  belong to labels of  $c$ . By putting suitable rules on how sets of states occur as labels of  $c$ , we can enforce  $r$  to respect equality constraints.

**Definition 6 (configurations).** *A configuration  $c$  of  $A$  is a rooted directed acyclic graph such that: (i) every node carries a symbol from  $2^Q$ , (ii) outgoing edges of a node are ordered, and (iii) for every equivalence class  $[q]$  of  $=_A$ , there is at most one set  $P \subseteq Q$  such that  $[q] \cap P \neq \emptyset$  and  $P$  is a label of  $c$ .*

Nodes of  $c$  are denoted by  $\text{nodes}(c)$ . For every node  $u \in \text{nodes}(c)$ , we denote by  $c|_u$  the subgraph of  $c$  induced by the nodes reachable from  $u$  in  $c$ , and by  $u_1, \dots, u_n$  the  $n$  successor nodes of  $u$  given in order. Note that it might be the case that  $u_i = u_j$  for some  $i, j \in \{1, \dots, n\}$ . Finally, we always denote by  $\text{lab}_c(u) \subseteq Q$  the label of node  $u$  in  $c$ . Fig. 3 illustrates a configuration whose nodes are labelled either by set of states  $\{q, q'\}$ ,  $\{s\}$ ,  $\{p\}$  or  $\{r\}$ . Note that the second successor of the root is the node labelled  $\{r\}$ , while its first and third successor is the node labelled  $\{q, q'\}$ .



**Fig. 3.** A configuration (naturals represent the order on edges)

In order to define semantics of configurations, we first introduce some useful notions about contexts. For  $n \geq 0$ , we define  $n$ -ary contexts  $C$ s as usual, and the hole substitution with trees  $t_1, \dots, t_n$  is denoted by  $C[t_1, \dots, t_n]$ . Note that 0-ary contexts are just trees. See [10] for a formal definition. Given  $n$  states  $p_1, \dots, p_n \in Q$  and an  $n$ -ary context  $C$ , we denote by  $C[p_1, \dots, p_n]$  the tree over  $\Lambda \cup Q$ , where each  $p_i$  is

viewed as a constant symbol. We let  $A^*$  be the TAGED over alphabet  $A \cup Q$  which is just  $A$  with additional rules  $q(\epsilon) \rightarrow q$ , for every  $q \in Q$ . We say that  $C[p_1, \dots, p_n]$  evaluates to  $p$ , denoted  $C[p_1, \dots, p_n] \rightarrow_A p$  if there is a run  $r$  of  $A^*$  on  $C[p_1, \dots, p_n]$  such that  $\text{roots}(r) = p$ . For any set  $S \subseteq Q$ , we write  $C[p_1, \dots, p_n] \rightarrow_{Q \setminus S} q$  if states from  $S$  does not occur in  $r$ , except at the root and at the leaves labelled  $p_1, \dots, p_n$ .

We now view configurations as a way to combine contexts to form trees  $t$ , with additional requirements which enforce existence of a run  $r$  of  $A$  on  $t$ . Condition (iii) of Definition 6 ensures  $r$  satisfies the equality constraints. This motivates the semantics of configurations. More formally, let  $c$  be a configuration of  $A$ . A mapping  $\lambda$  from  $\text{nodes}(c)$  into contexts over  $A$  is an *interpretation* of  $c$  if for every node  $u \in \text{nodes}(c)$ , if  $u$  has exactly  $n$  successor nodes  $u_1, \dots, u_n$  in  $c$ , then  $\lambda(u)$  is an  $n$ -ary context. Moreover,  $\lambda$  must satisfy the following: for every node  $u \in \text{nodes}(c)$  and every nodes  $u_1, \dots, u_n \in \text{nodes}(c)$ , if  $u_1, \dots, u_n$  are the successor nodes of  $u$  then the following holds (called condition (P)):

$$\forall p \in \text{lab}_c(u), \exists p_1 \in \text{lab}_c(u_1) \dots \exists p_n \in \text{lab}_c(u_n), \lambda(u)[p_1, \dots, p_n] \rightarrow_{Q \setminus (E \cup D)} p$$

As  $A$  is positive, the set  $D$  is empty, but we keep this definition for uniformity reasons when dealing with disequalities. Every node  $u \in \text{nodes}(c)$ , together with the mapping  $\lambda$ , define a tree  $t(u, \lambda)$  over  $A$  as follows:  $t(u, \lambda) = \lambda(u)[t(u_1, \lambda), \dots, t(u_n, \lambda)]$ , where  $n \in \mathbb{N}$  and  $u_1, \dots, u_n$  are the successor nodes of  $u$  in  $c$ . Note that this definition is well-founded since  $c$  is a DAG. As we will see, condition (P) ensures the existence of a run of  $A$  on  $t(u_0, \lambda)$ , where  $u_0$  is the root of  $c$ .

**Definition 7 (tree language recognized by a configuration).** Let  $c$  be a configuration of  $A$ . The tree language recognized by  $c$ , denoted  $L(A, c)$ , is defined by the set of trees  $t(u_0, \lambda)$ , where  $u_0$  is the root of  $c$ , and  $\lambda$  is an interpretation of  $c$ .

Trees accepted by configuration of Fig. 3 are necessarily of form  $C[C'[t], t', C'[t]]$ , for some contexts  $C, C', C''$  and trees  $t, t'$ . As already said, the constraints on the contexts and the configuration ensure the existence of a run on the trees of  $L(A, c)$  which satisfies the equality constraints. In particular, we can prove the following:

**Proposition 1.** Let  $c$  be a configuration of  $A$  such that  $L(A, c)$  is nonempty. Let  $t$  be a tree of  $L(A, c)$ , and  $u_0 \in \text{nodes}(c)$  be the root of  $c$ . For every  $p \in \text{lab}_c(u_0)$ , there is a run  $r \in R_{A,p}(t)$  which respects the equality constraints.

The converse holds too, and we can bound the size of configurations:

**Proposition 2.** Let  $t \in \mathbb{T}_A$  be a tree such that  $t \in L(A)$ . Then there is a configuration  $c$  of size at most  $|Q| \cdot k^{|Q|}$  such that  $t \in L(A, c)$ .

*Sketch of proof* We start from an accepting run  $r$  of  $A$  on  $t$  and define an equivalence relation on a subset of  $\text{nodes}(t)$ . Informally, two nodes  $u, v$  are equivalent if an equality test between  $t|_u$  and  $t|_v$  is performed in  $r$ . This is the case for instance when  $\text{lab}_r(u) =_A \text{lab}_r(v)$ . Each equivalence class will represent a node of  $c$ , to enforce equalities.  $\square$

Hence, we can finitely represent the language recognized by  $A$  as a computable set of configurations of  $A$ , as stated below:

**Corollary 1.** *Let  $A$  be a  $k$ -bounded positive TAGED. We let  $\mathcal{D}(A)$  be the set of configurations of  $A$  whose sizes are bounded by  $|Q| \cdot k^{|Q|}$ . The following holds:*

$$L(A) = \bigcup_{c \in \mathcal{D}(A)} L(A, c)$$

Moreover, we can decide emptiness of the language recognized by any configuration.

**Lemma 1.** *Given a configuration  $c$ , it is decidable to know whether  $L(A, c) = \emptyset$  holds.*

*Proof (Sketch).* For all nodes  $u, u_1, \dots, u_n$  s.t.  $u_1, \dots, u_n$  are the successors of  $u$ , it suffices to test whether there is an  $n$ -ary context  $C$  s.t. for all state  $p \in \text{lab}_c(u)$ , there are  $p_1 \in \text{lab}_c(u_1), \dots, p_n \in \text{lab}_c(u_n)$  s.t.  $C[p_1, \dots, p_n] \rightarrow_{Q \setminus E} p$ . We can represent the set of contexts  $C$  such that  $C[p_1, \dots, p_n] \rightarrow_{Q \setminus E} p$  by a tree automaton  $A_{(p_i)_i, p}$ . Then, it suffices to test emptiness of  $\bigcap_{p \in P} \bigcup_{(p_i)_i \in \prod_i \text{lab}_c(u_i)} L(A_{(p_i)_i, p})$ , which is decidable, since regular languages are closed by Boolean operations.

As a corollary of Lemma 1 and Corollary 1, we get the following:

**Proposition 3.** *Emptiness of positive bounded TAGED is decidable.*

### 4.3 Adding disequalities to positive bounded TAGED

In the previous section, we have shown that emptiness of positive bounded TAGED is decidable. In this section, we extend this result to full bounded TAGED.  $A$  always denotes a  $k$ -bounded TAGED. The definition of configurations of  $A$  remains unchanged, and the set  $\mathcal{D}(A)$  still denotes the set of configurations of  $A$  whose size is bounded by  $|Q| \cdot k^{|Q|}$ . We have the following inclusion:  $L(A) \subseteq \bigcup_{c \in \mathcal{D}(A)} L(A, c)$ , but the other one does not hold in general, since configurations do not require disequality constraints to be satisfied. We show how an additional test on configurations  $c$  allows us to decide whether  $L(A) \cap L(A, c)$  is empty, which will be sufficient to decide whether  $L(A)$  is empty. Informally, let  $c$  be a configuration and  $\lambda$  be an interpretation of  $c$ . We say that  $\lambda$  satisfies the disequalities of  $c$  if for all nodes  $u, v \in \text{nodes}(c)$ , if there are  $p \in \text{lab}_c(u)$  and  $q \in \text{lab}_c(v)$  such that  $p \neq_A q$ , then  $t(u, \lambda) \neq t(v, \lambda)$ .

We now relate the problem of finding such an interpretation to context disunification. For all nodes  $u \in \text{nodes}(c)$ , we let  $\text{cxt}_c(u)$  be the set of contexts satisfying condition (P) of the definition of interpretation. We define the notion of *partial interpretation*  $\beta$  of  $c$ , as a mapping from  $\text{nodes}(c)$  into contexts, such that it maps every node  $u$  such that  $\text{cxt}_c(u)$  is finite into a context of  $\text{cxt}_c(u)$ , and every other node  $v$  to a context  $@_v(\bullet, \dots, \bullet)$  with  $n$  holes (if  $v$  has  $n$  successors), where  $@_v$  is a fresh symbol such that  $@_v \notin \Lambda$ . Note that trees  $t(u, \beta)$  are trees over alphabet  $\Lambda \cup \{@_v \mid v \in \text{nodes}(c)\}$ . We can show the following, by using context disunification (symbols  $@_v$  are viewed as context variables):

**Lemma 2.** *Let  $A$  be a bounded TAGED. We have  $L(A) \neq \emptyset$  iff there exist a configuration  $c \in \mathcal{D}(A)$  and a partial interpretation  $\beta$  of  $c$  such that  $\beta$  satisfies the disequality constraints of  $c$ . Moreover, it is decidable to know whether such an interpretation  $\beta$  exists.*

As a corollary, by combining Lemma 2 and Lemma 1, we get the proof of Theorem 5.

## 5 From TQL to Automata

In this section,  $\phi$  denotes a recursion-closed and guarded TQL formula over tree variables  $X_1, \dots, X_n$ . We sketch the construction of a TAGED  $A_\phi$  such that for all hedges  $h \in \mathbb{H}_A$ , we have  $h \in L(A_\phi)$  iff there exists a valuation  $\rho : \text{var}(\phi) \rightarrow \mathbb{T}_A$  such that  $h \in \llbracket \phi \rrbracket_\rho$ . Moreover, we prove  $A_\phi$  to be equivalent to a computable bounded TAGED whenever  $\phi$  is bounded.

In a first step, we transform  $\phi$  into an equivalent system of fixpoint equations, and then sketch the construction of  $A_\phi$  starting from this system. This construction extends the construction of [4]. This extension is non-trivial, since it manages tree variables, which induce non-determinism in the produced tree automaton. Moreover, even for sentences, this construction is different, since trees are ordered.

*System of equations* We define dual connectives for parallel composition and Kleene star. We let  $\phi_1 || \phi_2$  as a shortcut for  $\neg(\neg\phi_1 | \neg\phi_2)$ ,  $\phi_1^\diamond$  for  $\neg(\neg\phi_1)^*$  and  $\overline{X}$  for  $\neg X \wedge \Lambda[\top]$ . A *system of fixpoint equations*  $\Sigma$  is a sequence of equations  $\xi_1 = \text{rhs}_1, \dots, \xi_n = \text{rhs}_n$  where every  $\text{rhs}_i$  has one of the following form:

$$0 \quad \overline{0} \quad \xi \vee \xi' \quad \xi \wedge \xi' \quad \alpha[\xi] \quad \xi | \xi' \quad \xi || \xi' \quad X \quad \overline{X} \quad \xi^* \quad \xi^\diamond$$

The last fixpoint variable,  $\xi_n$ , is denoted by  $\text{last}(\Sigma)$ . The set of tree variables occurring in  $\Sigma$  is denoted by  $\text{var}(\Sigma)$ . Systems of equations are interpreted over the complete lattice  $(2^{\mathbb{H}_A}, \cup, \cap)$ , modulo an assignment  $\rho : \text{var}(\Sigma) \rightarrow \mathbb{T}_A$ . We consider the following monotonic operations over  $2^{\mathbb{H}_A}$ , modulo  $\rho$ :  $0$  is interpreted as  $\{0^{\mathbb{H}_A}\}$ ,  $\overline{0}$  as  $\{\overline{0^{\mathbb{H}_A}}\}$  (the overline denotes the complement in  $\mathbb{H}_A$ ),  $\vee$  by  $\cup$ ,  $\wedge$  by  $\cap$ ,  $\alpha[\cdot]$  as the unary operator which maps any set of hedges  $H \subseteq \mathbb{H}_A$  into  $\{a[h] \mid a \in \alpha, h \in H\}$ ,  $|\cdot|$  as the binary operator which maps two sets of hedges  $H, H'$  into  $H | H' = \{h | h' \mid h \in H, h' \in H'\}$ , its dual  $||\cdot|$  maps  $H, H'$  into  $H || H' = \overline{H | H'}$ . The Kleene star  $^*$  and its dual  $^\diamond$  are interpreted similarly. Finally,  $X$  is interpreted by  $\rho(X)$  and  $\overline{X}$  by  $\mathbb{T}_A \setminus \{\rho(X)\}$ .

The solution of  $\Sigma$  over  $(2^{\mathbb{H}_A}, \cup, \cap)$  modulo  $\rho$  is a mapping from fixpoint variables of  $\Sigma$  into  $2^{\mathbb{H}_A}$ , and is denoted by  $\text{Sol}_{\mathbb{H}}(\Sigma, \rho)$ . We can push down the negations to the leaves of  $\phi$ , using the dual connectives, and introduce fixpoint variables for every position in  $\phi$ , which allow to construct a system of equations  $S_\phi$  such that  $\text{var}(S_\phi) = \text{var}(\phi)$  and the following holds:

**Lemma 3.** *For all valuations  $\rho : \text{var}(S_\phi) \rightarrow \mathbb{T}_A$ , we have  $\llbracket \phi \rrbracket_\rho = \text{Sol}_{\mathbb{H}}(S_\phi, \rho)(\text{last}(S_\phi))$*

E.g., the equation system associated with the formula  $\mu\xi.(a[\xi] \vee (\mu\xi'.(b[\xi'] \vee X)))$  is  $\{\xi' = \xi_2 \vee \xi_3; \quad \xi_2 = b[\xi']; \quad \xi_1 = a[\xi]; \quad \xi_3 = X; \quad \xi = \xi_1 \vee \xi'\}$ .

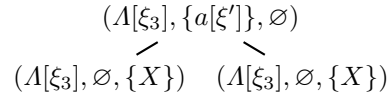
*Ideas of automaton construction for sentences* In this paragraph, we assume that  $\phi$  is a sentence. Checking whether a hedge  $h$  is a solution of  $S_\phi$  is similar to a run of some hedge automaton on  $h$ . Let us consider the system  $S = \{\xi = \xi_1 \vee \xi_2; \xi_1 = a[\xi]; \xi_2 = 0\}$ , where the last variable is  $\xi$ . Solutions of  $S$  are chains labelled by  $as$ . To check whether hedge  $a(0)$  is a solution of  $\xi$ , first verify that  $a(0)$  is a solution of  $\xi_1$  or a solution of  $\xi_2$ . One can easily see that  $a(0)$  is not a solution of  $\xi_2$ . It remains to verify whether  $a(0)$  is a solution of  $\xi_1 = a[\xi]$ , by verifying that  $0$  is a solution of  $\xi$ , etc... This can be done by an automaton with transition rules  $a(\epsilon + q_{a[\xi]}) \rightarrow q_{a[\xi]}$ , where  $q_{a[\xi]}$  is

a final state. We define the set of (final) states by  $Q = F = \{q_{a[\xi]}\}$ . Let us interpret  $S$  over  $2^{Q^*}$ , where  $Q^*$  is the set of words over  $Q$ . We interpret  $\vee$  by  $\cup$ ,  $0$  by  $\{\epsilon\}$ , and  $a[\xi]$  by  $\{q_{a[\xi]}\}$ . Solutions of  $\xi$  are denoted by  $s(\xi)$  (and similarly for other variables). Hence we get  $s(\xi_2) = \{\epsilon\}$ ,  $s(\xi_1) = \{q_{a[\xi]}\}$ , and  $s(\xi) = \{\epsilon, q_{a[\xi]}\}$ . Which trees evaluate to  $q_{a[\xi]}$ ? Trees of the form  $a(h)$  where  $h$  evaluates to some word of states from  $s(\xi)$ . Hence, we can define transition  $a(s(\xi)) \rightarrow q_{a[\xi]}$ .

Things get more complicated when adding intersection. For instance, consider the system  $S' = \{\xi_1 = a[\xi]; \xi_2 = a[\xi']; \xi' = 0; \xi = \xi_1 \wedge \xi_2\}$ . If we interpret this system as before, with states  $q_{a[\xi]}$  and  $q_{a[\xi']}$ , we would get  $s(\xi) = \emptyset$ . Hence, states should carry enough information to know if the current tree is a solution of  $a[\xi]$ ,  $a[\xi']$ , or both. In the construction we provide, every state is a tuple of atoms of the form  $\alpha[\xi]$ ,  $\overline{\alpha}[\top]$ , or  $\alpha[\neg\xi]$ , for every right-hand side of the form  $\alpha[\xi]$  occurring in the system. If some tree  $t$  evaluates to a tuple which has a component equal to  $\alpha[\xi]$ , it means that  $t$  is of the form  $a(h)$ , where  $a \in \alpha$  and  $h$  is solution of  $\xi$ . If the component is  $\overline{\alpha}[\top]$  or  $\alpha[\neg\xi]$ , it means that  $a(h)$  is not a solution of  $\alpha[\xi]$ , because, in the first case,  $a \notin \alpha$ , and in the second one,  $a \in \alpha$ , but  $h$  is not a solution of  $\xi$ . Knowing this complete information, i.e. which right-hand sides of the form  $\alpha[\xi]$  are satisfied or not by the current tree, we are able to construct exactly one rule per state, by solving the system on sets of words of states, with suitable interpretations. As the formula is guarded, solutions of the system are regular state word languages. We then get a deterministic hedge automaton whose accepted trees are the solutions of the system.

*Adding tree variables* When adding variables, we cannot keep the automaton deterministic, since subtrees will be captured non-deterministically. For instance, consider the system  $S = \{\xi'' = X; \xi' = \xi''|\xi''; \xi_X = a[\xi']; \xi_2 = 0; \xi_3 = \xi_+; \xi_1 = \Lambda[\xi_3]; \xi_\top = \xi_1 \vee \xi_2; \xi = \xi_\top \wedge \xi_X\}$ , where the last variable is  $\xi$ . Given a valuation  $\rho : \text{var}(S) \rightarrow \mathbb{T}_A$ , the system  $S$  has a unique solution, modulo  $\rho$ , which is  $a(\rho(X)|\rho(X))$ . A TAGED accepting the solutions of  $S$  is the TAGED of Example 1 of Section 4. It is non-deterministic, since it has to choose to go in a state  $q_X$  which will enforce the TAGED to test whether the two sons of the root are equal.

As tree variables induce a kind of non-determinism, we emphasize two kinds of recursion variables: deterministic recursion variables, for which the problem of checking whether a given hedge is a solution does not involve tree variables, such as  $\xi_\top$ ,  $\xi_1$ ,  $\xi_2$  and  $\xi_3$  in our example; and non-deterministic one:  $\xi, \xi_X, \xi'$  and  $\xi''$ .



**Fig. 4.** Run of  $A_\phi$  on  $a(a(0)a(0))$

States of the automaton we construct have three components. The first component is induced by deterministic recursion variables and simulate a classical hedge automaton, as for the case of sentences. The second component is induced by non-deterministic recursion variables, and collects atoms of the form  $\alpha[\xi]$ , where  $\xi$  is non-deterministic, for which the current tree is the solution. In other words, it guesses the positions in the tree under which capture variables occur. Third components are sets of variables  $X$  or  $\overline{X}$ , meaning that equality or disequality tests are performed on the current tree.

Transition rules are obtained by suitable interpretation of the system over words of states. Finally, two states are equivalent for the automaton if their third component

shares a tree variable. Disequalities are defined similarly. For instance, an accepting run of the automaton for  $S$ , on the tree  $a(a(0)a(0))$ , is represented in Fig. 4. The state  $(A[\xi_3], \emptyset, \{X\})$  is equivalent to itself.

**Theorem 6.** *Let  $\phi$  be a guarded TQL formula. There is a computable TAGED  $A_\phi$  such that for all hedges  $h$ , we have  $h \in L(A_\phi)$  iff there exists a valuation  $\rho : \text{var}(\phi) \rightarrow \mathbb{T}_A$  such that  $h \in \llbracket \phi \rrbracket_\rho$ .*

*Moreover, if  $\phi$  is bounded, the TAGED  $A_\phi$  is equivalent to the bounded TAGED  $(A_\phi, B)$ , for some computable bound  $B \in \mathbb{N}$ .*

To compute the bound  $B$ , we interpret the system  $S_\phi$  on naturals, with suitable interpretations (for instance,  $X$  is interpreted by 1,  $\wedge$  by  $+$ , and  $\vee$  by  $\max$ ).

## 6 Extending MSO with Tree Isomorphism Tests

In this section, we propose an extension of MSO for unranked trees with isomorphism tests between trees.

Let  $\sim$  be a binary predicate s.t. for a structure  $\mathcal{S}^h$  associated with a hedge  $h$  and a mapping  $\rho$  from  $\{x, y\}$  to nodes of  $h$ ,  $\mathcal{S}^h \models_\rho x \sim y$  holds if the two subtrees rooted at respectively  $\rho(x)$  and  $\rho(y)$  in  $h$  are isomorphic. We consider sentences of the form

$$Q_1 x_1 Q_2 x_2 \dots Q_n x_n \psi$$

where  $Q_i \in \{\exists, \forall\}$  and  $\psi$  is an MSO formula extended with atoms  $x_i \sim x_j$  ( $1 \leq i, j \leq n$ ). We call  $MSO_\sim$  this logic. We will also consider the fragment  $MSO_\sim^\exists$  for which formulas satisfy  $Q_1 = Q_2 = \dots = Q_n = \exists$ . Remark that  $\psi$  can again contain quantifiers. Hence, since tree isomorphism cannot be expressed in  $MSO$  [10],  $MSO_\sim$  and  $MSO_\sim^\exists$  are strictly more expressive than  $MSO$ . By reduction of the Post correspondence problem, we can prove that:

**Theorem 7.** *Satisfiability for  $MSO_\sim$  is undecidable.*

However,  $MSO_\sim^\exists$  and bounded TAGED are equally expressive: for any formula  $\varphi$  in  $MSO_\sim^\exists$ , one can compute a bounded TAGED, whose size is non-elementary in the size of  $\varphi$ , accepting the models of  $\varphi$ . The converse holds too. As a consequence of decidability of emptiness for bounded TAGED, we have:

**Theorem 8.** *Satisfiability is decidable for  $MSO_\sim^\exists$ .*

## 7 Conclusion

In this paper, we have considered the spatial logic TQL with tree variables. We have proved that for the guarded fragment when variables appear in a bounded way then the satisfiability problem is decidable. To do so, we have introduced a new class of tree automata, called TAGED, permitting to test global equalities and disequalities on the accepted trees. Finally, we have used TAGED to prove decidability for an extension of MSO with isomorphism tests interpreted over unranked trees.

We speculate that the boundedness condition is not required for the decidability of emptiness of TAGED, as pumping techniques dealing with constraints may be applicable. However, it is non-trivial, since TAGED are not determinizable in general. This would imply that the full guarded TQL with trees variables is decidable. Another extension

would be to consider hedge variables. This problem seems to be non trivial as the satisfiability problem for such formulas could encode word equations.

We emphasized a correspondence between  $MSO_{\sim}^{\exists}$  and bounded TAGED. It would be interesting to exhibit a fragment of  $MSO_{\sim}$  equivalent to full TAGED.

The TQL system also includes a transformation language; we aim at using TAGED automata to type these transformations and more generally, tree transducers.

## References

1. V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. In *8th ACM International Conf. on Functional Programming*, pages 51–63, 2003.
2. B. Bogaert and S. Tison. Equality and disequality constraints on direct subterms in tree automata. In *9th Annual Symposium on Theoretical Aspects of Computer Science*, volume 577 of *LNCS*, pages 161–171, 1992.
3. M. Bojanczyk, C. David, A. Muscholl, T. Schwentick, and L. Segoufin. Two-variable logic on data trees and xml reasoning. In *ACM 25th Symp. on Principles of database systems*, pages 10–19, 2006.
4. I. Boneva, JM. Talbot, and S. Tison. Expressiveness of a spatial logic for trees. In *20th IEEE Symposium on Logic in Computer Science*, pages 280–289, 2005.
5. A. Brüggemann-Klein, M. Murata, and D. Wood. Regular tree languages over non-ranked alphabets. unpublished manuscript, 1998.
6. L. Cardelli, P. Gardner, and G. Ghelli. A spatial logic for querying graphs. In *29th International Colloquium on Automata, Languages and Programming*, volume 2380 of *LNCS*, pages 597–610. Springer, 2002.
7. L. Cardelli and G. Ghelli. TQL: A Query Language for Semistructured Data Based on the Ambient Logic. *Mathematical Structures in Computer Science*, 14:285–327, 2004.
8. W. Charatonik and JM. Talbot. The Decidability of Model Checking Mobile Ambients. In *15th Annual Conference of the European Association for Computer Science Logic*, volume 2142 of *LNCS*, pages 339–354. Springer, 2001.
9. H. Comon and V. Cortier. Tree automata with one memory, set constraints and cryptographic protocols. *TCS*, 331(1):143–214, 2005.
10. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications, 1997. release October, 1st 2002.
11. M. Dauchet, A.-C. Caron, and J.-L. Coquidé. Reduction properties and automata with constraints. 20:215–233, 1995.
12. A. Dawar, P. Gardner, and G. Ghelli. Expressiveness and complexity of graph logic. In *Information and Computation*, number 205, pages 263–310. 2007.
13. H. Hosoya and B. C. Pierce. XDuce: A statically typed xml processing language. *ACM Trans. Internet Techn.*, 3(2):117–148, 2003.
14. F. Jacquemard, M. Rusinowitch, and L. Vigneron. Tree automata with equality constraints modulo equational theories. Research Report LSV-06-07, LSV, ENS Cachan, France, 2006.
15. W. Karianto and C. Löding. Unranked tree automata with sibling equalities and disequalities. In *34th International Colloquium on Automata, Languages and Programming*, 2007.
16. J. Mongy. *Transformation de noyaux reconnaissables d'arbres. Forêts RATEG*. PhD thesis, Université de Lille, 1981.
17. M. Murata. Hedge automata: A formal model for xml schemata. Technical report, Fuji Xerox Information Systems, 1999.
18. Frank Neven, Thomas Schwentick, and Victor Vianu. Towards regular languages over infinite alphabets. In *MFCS*, pages 560–572. SV, 2001.
19. J. C. Reynolds. Separation logic: A logic for shared mutable data structures. In *17th IEEE Symp. on Logic in Computer Science*, pages 55–74. IEEE, 2002.